

Зачётный тест по курсу «Байесовский выбор моделей»

Время выполнения: 90 минут

Максимальный балл: 100 баллов, поэтому можно выполнять не все задания

Замечание: попавшие в TOP2 по баллам дополнительно получают $\min(X, 100 - 50R)$, где X – набранный балл, а R – ранг от 0 до 1.

Задача 1 (20 баллов). Пусть имеется НОР (i.i.d.) выборка x_1, \dots, x_n из нормального распределения $\mathcal{N}(m, \sigma^2)$, то есть $x_i \sim \mathcal{N}(x_i|m, \sigma^2)$. Введем априорные распределения на m и σ^2 вида

$$m \sim \mathcal{N}(m|m_0, \sigma_0^2), \frac{1}{\sigma^2} \sim \Gamma\left(\frac{1}{\sigma^2}|\alpha, \beta\right), \quad (1)$$

где $\alpha, \beta, m_0, \sigma_0^2$ – известные гиперпараметры.

а) Выписать совместное правдоподобие модели $p(\mathbf{x}, m, \sigma^2|\alpha, \beta, m_0, \sigma_0)$ (1 балл);

б) Выписать $p(m, \sigma^2|\mathbf{x}, \alpha, \beta, m_0, \sigma_0)$. Принадлежит ли оно параметрическому семейству нормально-гамма распределений (normal-gamma distributions) и почему? (3 балла);

в) Является ли исходное априорное распределение (1) сопряженным к правдоподобию и почему? (1 балл) Найти и обосновать параметрическое семейство сопряженных распределений (2 балла);

г) Получить вариационное приближение $q(m, \sigma^2) = q(m)q(\sigma^2)$ для полного апостериорного распределения $p(m, \sigma^2|\mathbf{x}, \alpha, \beta, m_0, \sigma_0)$ (13 баллов).

Задача 2 (40 баллов). Пусть имеется K математических кружков, в каждом из которых обучается N студентов. Пусть известны результаты решения одинаковых по сложности задач студентами кружков. Там для студента с номером n в кружке с номером k известны два числа: t_{kn} – количество задач, которые студент попробовал решить и s_{kn} – количество успешно решенных задач. В каждом кружке считаем, что вероятность каждого студента решить задачу равна p_k , зависящая от кружка, но не зависящая от студента. Успешность решения разных задач одним студентом, а также успешность решения задач между студентами независимы в совокупности. Считаем, что априорное распределение на p_k есть $p_k \sim \text{Beta}(p_k|\alpha, \beta)$, где α, β – неизвестные гиперпараметры.

Обозначим $\mathbf{p} = [p_k, k = 1, \dots, K]^T$, $\mathbf{T} = \|t_{kn}\|$, $\mathbf{S} = \|s_{kn}\|$, $k = 1, \dots, K$, $n = 1, \dots, N$.

а) Выписать в явном виде совместное правдоподобие $p(\mathbf{S}, \mathbf{p}|\mathbf{T}, \alpha, \beta)$ (5 баллов);

б) Получить апостериорное распределение $p(\mathbf{p}|\mathbf{T}, \mathbf{S}, \alpha, \beta)$ и описать структуру зависимостей между компонентами \mathbf{p} (10 баллов). Вычислить $\mathbb{E}\mathbf{p}$ по апостериорному распределению (3 балла). Какие выводы можно сделать из полученного результата? (2 балла)

в) Выписать обоснованность $p(\mathbf{S}, \mathbf{T}|\alpha, \beta)$ в явном виде. Описать, как найти оценки гиперпараметров α, β из принципа максимума обоснованности (20 баллов).

Задача 3 (20 баллов). Пусть рассматривается поток посетителей в магазин, и измеряются интервалы между приходом двух последовательных посетителей. Считаем, что интервалы t_1, \dots, t_k, \dots между последовательными посетителями независимы в совокупности имеют показательное распределение с параметром λ , то есть

$$p(t_k) = \lambda \exp(-\lambda t_k), t_k \geq 0.$$

а) Выписать правдоподобие модели $p(\mathbf{t}|\lambda)$ (2 балла). Ввести априорное распределение на λ $p(\lambda|\alpha)$, сопряженное с правдоподобием, где α – вектор гиперпараметров (3 балла). Какое семейство распределений сопряжено с таким правдоподобием? (1 балл)

б) Получить выражения для обоснованности модели $p(\mathbf{t}|\alpha)$ в явном виде и описать метод поиска α из принципа максимума обоснованности (14 баллов).

Задача 4 (40 баллов). Пусть имеется обучающая выборка (\mathbf{X}, \mathbf{y}) , $\mathbf{X} \in \mathbb{R}^{m \times n}$, $\mathbf{y} \in$

$[-1, 1]^m$, полученная из модели генерации данных с совместным правдоподобием

$$p(\mathbf{y}, \mathbf{w}|\mathbf{A}, \mathbf{X}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{A}^{-1}) \prod_j p(y_j|\mathbf{x}_j, \mathbf{w}),$$

где $p(y_j|\mathbf{x}_j, \mathbf{w})$ дается моделью логистической регрессии, то есть

$$\mathbb{P}(y_j = 1) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x}_j)} = \sigma(\mathbf{w}^\top \mathbf{x}_j), \quad \mathbf{A} = \text{diag}(\alpha_j).$$

Используем принцип максимума обоснованности для отбора признаков

$$\mathbf{A} = \arg \max_{\mathbf{A}} p(\mathbf{y}_1|\mathbf{X}_1, \mathbf{A}).$$

- Выписать совместное правдоподобие модели $p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \mathbf{A})$ в явном виде (1 балл);
- Какое значение α_j соответствует тому, что признак j незначим, и не используется в модели? (1 балл)
- Использовать вариационную нижнюю оценку для сигмоидной функции (см. лекцию 7) для получения нижней оценки на совместное правдоподобие (требуется привести подробный вывод) $p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \mathbf{A}) \geq L(\mathbf{w}, \mathbf{A}, \boldsymbol{\xi})$, где $\boldsymbol{\xi} \in \mathbb{R}^m$ – вектор дополнительных переменных (8 баллов);
- Для нижней оценки на обоснованность

$$p(\mathbf{y}|\mathbf{X}, \mathbf{A}) = \int p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \mathbf{A}) d\mathbf{w} \geq \int L(\mathbf{w}, \mathbf{A}, \boldsymbol{\xi}) d\mathbf{w} = \tilde{L}(\mathbf{A}, \boldsymbol{\xi})$$

получить формулы EM-алгоритма для решения задачи ее максимизации

$$\tilde{L}(\mathbf{A}, \boldsymbol{\xi}) \rightarrow \max_{\mathbf{A}, \boldsymbol{\xi}} \quad (30 \text{ баллов}).$$

Задача 5 (70 баллов). Пусть рассматривается поток посетителей в магазин, и измеряются интервалы между приходом двух последовательных посетителей. Считаем, что распределение интервала t_n между последовательными посетителями описывается смесью K показательных распределений с показателями $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_K]^\top$ и весами $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]^\top$, то есть

$$p(t_n) = \sum_{k=1}^K \pi_k \lambda_k \exp(-\lambda_k t_n),$$

причем наблюдаемые интервалы $\mathbf{t} = [t_1, \dots, t_N]^\top$ независимы в совокупности.

Введем на $\boldsymbol{\pi}$ априорное симметричное распределение Дирихле $\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\pi}|\boldsymbol{\mu}\mathbf{e})$, где $\mu < 1$ – параметр распределения Дирихле (считается фиксированным и известным), а \mathbf{e} – единичный вектор.

Введем также априорные распределения на λ_k вида

$$\lambda_k \sim \Gamma(\alpha, \beta), \quad k = 1, \dots, K,$$

где α, β – неизвестные параметры гамма-распределения.

- Выписать совместное правдоподобие модели $p(\mathbf{t}, \boldsymbol{\pi}, \boldsymbol{\lambda}|\alpha, \beta, \mu)$ (5 баллов);
- Ввести матрицу скрытых переменных \mathbf{Z} принадлежности объекта компоненте смеси и выписать совместное правдоподобие модели со скрытой переменной $p(\mathbf{t}, \boldsymbol{\pi}, \boldsymbol{\lambda}, \mathbf{Z}|\alpha, \beta, \mu)$ (5 баллов);
- Воспользовавшись вариационным EM-алгоритмом, получить аппроксимацию $q(\boldsymbol{\pi}, \boldsymbol{\lambda}, \mathbf{Z}) = q(\boldsymbol{\pi})q(\boldsymbol{\lambda})q(\mathbf{Z})$ для истинного апостериорного распределения $p(\boldsymbol{\pi}, \boldsymbol{\lambda}, \mathbf{Z}|\mathbf{t}, \mu, \alpha, \beta)$ (35 баллов), а также задачи для нахождения оценок максимума обоснованности для α, β

$$\alpha^*, \beta^* = \arg \max_{\alpha, \beta} p(\mathbf{t}|\mu, \alpha, \beta)$$

и описать способ решения полученных оптимизационных задач (25 баллов).